



# Nonlinear dimensionality reduction in climate data

A. J. Gámez, C. S. Zhou, A. Timmermann, J. Kurths

## ► To cite this version:

A. J. Gámez, C. S. Zhou, A. Timmermann, J. Kurths. Nonlinear dimensionality reduction in climate data. *Nonlinear Processes in Geophysics*, 2004, 11 (3), pp.393-398. hal-00302357

**HAL Id: hal-00302357**

**<https://hal.science/hal-00302357>**

Submitted on 13 Sep 2004

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Nonlinear dimensionality reduction in climate data

A. J. Gámez<sup>1</sup>, C. S. Zhou<sup>1</sup>, A. Timmermann<sup>2</sup>, and J. Kurths<sup>1</sup>

<sup>1</sup>Institut für Physik, Universität Potsdam, Postfach 601553, D-14415 Potsdam, Germany

<sup>2</sup>Leibniz Institut für Meereswissenschaften, IfM-GEOMAR, Düsterbrook Weg 20, D-24105 Kiel, Germany

Received: 20 July 2004 – Accepted: 18 August 2004 – Published: 13 September 2004

Part of Special Issue “Nonlinear analysis of multivariate geoscientific data – advanced methods, theory and application”

**Abstract.** Linear methods of dimensionality reduction are useful tools for handling and interpreting high dimensional data. However, the cumulative variance explained by each of the subspaces in which the data space is decomposed may show a slow convergence that makes the selection of a proper minimum number of subspaces for successfully representing the variability of the process ambiguous. The use of nonlinear methods can improve the embedding of multivariate data into lower dimensional manifolds. In this article, a nonlinear method for dimensionality reduction, Isomap, is applied to the sea surface temperature and thermocline data in the tropical Pacific Ocean, where the El Niño-Southern Oscillation (ENSO) phenomenon and the annual cycle phenomena interact. Isomap gives a more accurate description of the manifold dimensionality of the physical system. The knowledge of the minimum number of dimensions is expected to improve the development of low dimensional models for understanding and predicting ENSO.

## 1 Introduction

The reduction of dimensionality of large multivariate data is a common task in many different branches of science. Its purpose is to reduce a theoretically infinite dimensional system (which, in measured data, is finite but highly dimensional) to a few physically relevant modes in order to gain insight into the dynamics of the complex system by neglecting unimportant degrees of freedom. In climate research, a widely used linear method is the principal component analysis (PCA) (Jolliffe, 1986), also called empirical orthogonal function (EOF) analysis in the geoscience context (von Storch and Zwiers, 1999), proper orthogonal decomposition (POD) in fluid dynamics (Holmes et al., 1997) and Karhunen-Loève decomposition in the continuous form (Karhunen, 1946; Loève, 1945). In PCA, the system under study is approximated by

a linear combination of steady spatial patterns with time dependent coefficients. The relevant number of dimensions is determined from the cumulative explained variance of the PCA modes and their physical interpretation. This variance is closely related to the eigenvalues of the spectral decomposition. Unfortunately, for many physical systems, these eigenvalues show a slow convergence that hampers the selection of the minimum number of dimensions.

In this article, we will apply a nonlinear method of dimensionality reduction to the observed sea surface temperature (SST) data in the tropical Pacific Ocean. In the equatorial Pacific, the SST evolution is characterized by a nonlinear superposition of two different oscillatory phenomena, ENSO and the annual cycle. ENSO is the most dominant statistical and physical mode of climate variability on interannual timescales (Philander, 1990). Climate models of different complexity have been used to explore the origin of its oscillatory character, its period and skewness (Zebiak and Cane, 1987; Tziperman et al., 1994; Jin, 1997). Attempts to reconstruct ENSO's attractor have been made using different nonlinear methods (Monahan, 2001; Grieger and Latif, 1994). The statistical analysis of ENSO is mostly based on SST anomalies which are obtained by subtracting a mean annual cycle from the monthly averaged SST data. The annual cycle in the tropical Pacific area originates from a complex interplay between semi-annual solar forcing and coupled air-sea instabilities (e.g. Li and Philander, 1996; Xie, 1994). As the strength of these instabilities varies slowly in time, one may expect that the amplitude of the physical annual cycle is not stationary but time-dependent. Extracting a time-varying annual cycle and an ENSO mode in a multivariate way from SST data is not simple and the results may depend strongly on the assumptions used by different methodologies. In particular, linear methods may fail to disentangle both modes since ENSO and the annual cycle exhibit in some sense a joint synchronised behaviour – ENSO amplitude is strong during the boreal winter season. This behaviour is reminiscent of an interactive coupling between the two modes (Pikovsky et al., 2001). For this reason, the study of the inter-

Correspondence to: A. J. Gámez  
(gamez@agnld.uni-potsdam.de)

action is of great importance for understanding the variability of ENSO. In the last years, several articles discussing how ENSO and the annual cycle interact in the tropical Pacific Ocean have been published (e.g. Xie, 1995; Jin et al., 1996; Tziperman et al., 1998; An and Wang, 2001). This type of interaction, which could be nonlinear, may lead to erroneous conclusions when subtracting a constant annual cycle from SST data under consideration, as usually done in the analysis of ENSO dynamics. Therefore, the data space cannot be decomposed into a sum of linear subspaces each containing an independent variable because of the existing interaction. So that, the separation of the SST into physically independent modes is not possible. Our aim will be, then, to extract a low dimensional manifold where the whole physical system could be embedded.

The structure of the paper is as follows. In the second section, we will provide a view of the problem of dimensionality reduction using the framework of Multidimensional Scaling (MDS). By doing this, the linear (PCA) and nonlinear (Isomap) methods used in this article can be compared under the same theoretical basis. In fact, the crucial point will lie on the definition of distance in the data space. In the third section, the methods will be applied to the SST dataset. Finally, we will extract some conclusions about the feasibility of our procedure.

## 2 The framework of Multidimensional Scaling

Let us define a matrix  $\mathbf{T}_{n \times m}$  of data, in such a way that  $T_{ij}$  is the SST at time  $t_i$  in a certain spatial point  $\mathbf{x}_j$ . We can think that the matrix is built up of  $n$  points  $\{\mathbf{T}_i\}$  in a  $\mathbb{R}^m$  space, and that these points belong to a trajectory as they are ordered in time. In this space, for every two points  $\{\mathbf{T}_i, \mathbf{T}_j\}$  in  $\mathbb{R}^m$ , the metric or distance  $d_{ij} = d(\mathbf{T}_i, \mathbf{T}_j)$  can be defined as a function onto nonnegative real numbers that obey the following rules:  $d_{ij} = 0$  iff  $i = j$ ,  $d_{ij} = d_{ji}$  and the triangle inequality, which states that  $d_{ij} + d_{jk} \geq d_{ik}$ .

The distance can be interpreted as a measure of similarity: if  $d_{ij} < d_{ik}$  we say that the physical state represented by  $\mathbf{T}_i$  is closer to  $\mathbf{T}_j$  than  $\mathbf{T}_k$ . The multidimensional scaling (MDS) approach makes use of the matrix of squared distances

$$(\mathbf{D}^{(2)})_{ij} = d_{ij}^2 \quad (1)$$

and applies the following procedure (Borg and Groenen, 1997). First,  $\mathbf{D}^{(2)}$  is transformed via an operation called double-centering. More precisely, using a matrix  $\mathbf{J}$  where  $J_{ij} = \delta_{ij} - 1/n$ , where  $\delta_{ij}$  is the Kronecker delta, we define the matrix

$$\mathbf{Z}^{(2)} = -\frac{1}{2} \mathbf{J} \mathbf{D}^{(2)} \mathbf{J}. \quad (2)$$

As we will see later, the matrix  $\mathbf{Z}^{(2)}$  is the matrix of scalar products of data points in the case that we use euclidean distances. Then,  $\mathbf{Z}^{(2)}$  is decomposed into its eigenvalues and eigenvectors  $\mathbf{Z}^{(2)} \mathbf{p}_i = \lambda_i \mathbf{p}_i$ . Defining  $\Lambda_{ij} = \lambda_i \delta_{ij}$  with  $\lambda_i > \lambda_{i+1}$  and  $\mathbf{P}_{ij} = p_{ij}$  the component  $j$  of vector  $\mathbf{p}_i$ , then  $\mathbf{Z}^{(2)} = \mathbf{P} \Lambda \mathbf{P}^T$ . If the set of eigenvalues have some leading

components  $\{\lambda_1, \dots, \lambda_p\}$  with  $p < m$ , and the rest decays very fast, meaning  $\lambda_p \gg \lambda_{p+1}$ , we could approximate  $\mathbf{Z}^{(2)}$  to its projection into the subspace spanned by the eigenvectors  $\{\mathbf{p}_1, \dots, \mathbf{p}_p\}$ .

This MDS framework can be used to embrace all the data analysis methods we are going to use. In fact, the main difference between the linear and the nonlinear method applied in this work come, as we will show, from the particular definition of distance used in Eq. (1).

### 2.1 Principal Component Analysis

Let us consider our data space as euclidean. This means that, if we define the scalar product as,

$$\mathbf{T}_i \cdot \mathbf{T}_j = \sum_{k=1}^m T_{ik} T_{jk} \quad (3)$$

then, this definition allows us to define the norm of a vector as

$$\|\mathbf{T}_i\| = \sqrt{\mathbf{T}_i \cdot \mathbf{T}_i} = \sqrt{\sum_{k=1}^m T_{ik} T_{ik}} \quad (4)$$

and the distance between two points  $\mathbf{T}_i$  and  $\mathbf{T}_j$  as the norm of their difference

$$d_{ij} = \|\mathbf{T}_i - \mathbf{T}_j\| = \sqrt{\sum_{k=1}^m (T_{ik} - T_{jk})^2} \quad (5)$$

which is the well known expression of euclidean distance. Therefore,  $\mathbf{D}^{(2)}$  can be written in the form

$$\mathbf{D}^{(2)} = \mathbf{I} \mathbf{c}^T + \mathbf{c} \mathbf{I}^T - 2 \mathbf{T} \mathbf{T}^T \quad (6)$$

where the components of the vector  $\mathbf{c}$  are  $c_i = \sum_{k=1}^m T_{ik}^2$  and  $\mathbf{I}$  is a vector formed by ones and dimension  $n$ . Then, applying the centering and eigendecomposition as explained in Sect. 2, we obtain,

$$\mathbf{Z}^{(2)} = -\frac{1}{2} \mathbf{J} \mathbf{D}^{(2)} \mathbf{J} = \mathbf{T} \mathbf{T}^T = \mathbf{P} \Lambda \mathbf{P}^T = \mathbf{P} \Lambda^{1/2} \Lambda^{1/2} \mathbf{P}^T \quad (7)$$

where  $\Lambda$  is the diagonal matrix of eigenvalues and  $\mathbf{P}$  is the matrix of eigenvectors of  $\mathbf{Z}^{(2)}$ . This allows us to make a representation of  $\mathbf{T}$  in terms of  $\mathbf{P}_i \lambda_i^{1/2}$  which is called principal coordinates analysis (PCO or PCoA) (Gower, 1966). If the data have mean zero and variance one, there is a nice correspondence with the PCA. The correlation matrix  $\mathbf{V}$  can be written as,

$$\mathbf{V} = \mathbf{T} \mathbf{T} = \mathbf{Q} \Lambda \mathbf{Q}^T = \mathbf{Q} \Lambda^{1/2} \mathbf{P}^T \mathbf{P} \Lambda^{1/2} \mathbf{Q}^T \quad (8)$$

where  $\mathbf{P}^T \mathbf{P} = \mathbf{I}$ . Finally, we can associate, by singular value decomposition,  $\mathbf{T} = \mathbf{P} \Lambda^{1/2} \mathbf{Q}^T$ , where  $\mathbf{P} \Lambda^{1/2}$  is the principal components matrix,  $\Lambda$  is the diagonal matrix formed by the eigenvalues of the correlation matrix and  $\mathbf{Q}$  is the matrix of patterns. Therefore, PCA and PCO get the same representation in terms of components  $\mathbf{P} \Lambda^{1/2}$ . This way we can reduce the dimensionality of the data by taking  $\mathbf{P}_i \lambda_i^{1/2}$  from  $i = 1$  to some value  $p < m$ , which could be  $p \ll m$ .

Summarizing, PCA can be regarded as an euclidean MDS for normalised data. We would like to stress that PCA is a linear method of decomposition, where the data are projected into orthonormal linear subspaces. However, if the data points belong to a nonlinear manifold, the orthonormal projection spreads contributions to the variance onto the different principal components. In that case, reducing the dimensionality of a physical system in which the dynamics is not governed by linear processes or where there are nonlinear relations between variables could lead to a wrong interpretation of the dimensionality. As a simple example, we can apply PCA to several, but simple, nonlinear data. For instance, we construct a trajectory of a particle moving on a spiral with some added uniform noise (Fig. 1a). In this example we can observe (Fig. 1b) that the PCA do not reconstruct the one-dimensional trajectory. To reproduce the real dynamics we need, in this context, to use a definition of distance that captures the nonlinear structure of the manifold. The geodesic distance is the proper metric for measuring distances on nonlinear manifolds, as we will discuss in the following section.

## 2.2 Isomap

The natural metric for nonlinear manifolds is the geodesic distance (for a review in differential geometry, see Do Carmo (1976)). Let us have an euclidean space of temperatures  $\mathbf{T}$  of coordinates  $\{T_1, \dots, T_m\}$ . Suppose there is a manifold  $\Theta \subseteq \mathbf{T}$  represented by the coordinates  $\{\Theta_1, \dots, \Theta_p\}$ . The metric of  $\Theta$  is the matrix  $\mathbf{g}$  with elements defined by,

$$g_{ij} = \sum_{k=1}^m \frac{\partial T_k}{\partial \Theta_i} \frac{\partial T_k}{\partial \Theta_j} \quad (9)$$

For a general  $\mathbf{g}$ , the distance between two points  $\theta_1$  and  $\theta_2$  in  $\Theta$  is then

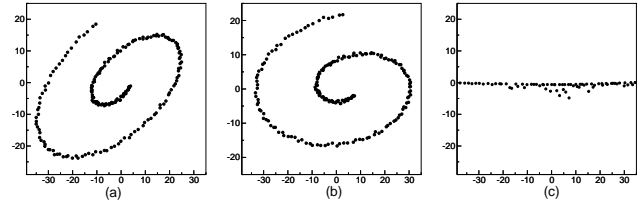
$$d(\theta_1, \theta_2) = \|\theta_1 - \theta_2\| = \int_{\theta_1}^{\theta_2} \sqrt{\sum_{i,j=1}^m g_{ij} d\Theta_i d\Theta_j} \quad (10)$$

In the case of a euclidean manifold,  $g_{ij} = \delta_{ij}$  and we recover Eq. (5) for the discrete case.

For example, let us think of a particular manifold  $\Theta$ , in our case a spiral similar to that one shown in Fig. 1a. The equations of the spiral are:

$$\begin{aligned} x &= t \sin t \\ y &= t \cos t, \end{aligned}$$

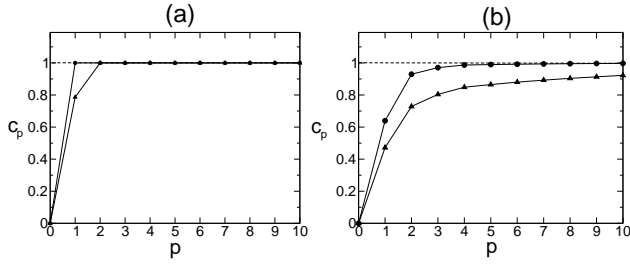
where  $t$  belongs to a real interval. The euclidean distance  $d_E$  between two points  $(x_1, y_1) = (x(t_1), y(t_1))$  and  $(x_2, y_2) = (x(t_2), y(t_2))$  is  $d_E = \sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2} = \sqrt{t_1^2 + t_2^2 - 2t_1t_2 \cos(t_1 - t_2)}$ . The path from one point to another one is a straight line that does not belong to the spiral. But, if we restrict ourselves to a path inside the spiral, the geodesic distance is  $d_S = \int_{t_1}^{t_2} \sqrt{1 + t^2} dt = \frac{1}{2} \{t_2 \sqrt{1 + t_2^2} - t_1 \sqrt{1 + t_1^2} + \operatorname{arcsinh}(t_2) - \operatorname{arcsinh}(t_1)\}$  as can be calculated. It is interesting to note that small euclidean distances may



**Fig. 1.** (a) Spiral with uniform noise. (b) Result for the spiral after PCA, the method cannot guess the true dimensionality. (c) Result for the spiral after Isomap, now the method unfolds the one dimensional trajectory.

correspond to large geodesic distances. For this reason, measuring similarities based on inadequate distances could lead to misleading results. Consequently, the idea is to substitute the euclidean distance in the MDS method (hence, the scalar product defined by the statistical correlation in the case of PCA) by the geodesic distance between each pair of points. Tenenbaum et al. (2000) proposed a method for computing geodesic distances through graph distances. The method could be divided in several phases: i) in a first step, Isomap approximates the geodesic distance using a graph constructed by connecting nearest neighbours in the euclidean  $\mathbb{R}^m$  space. More specifically, we will say that a point  $T_i$  is one of the nearest neighbours of  $T_j$  if it belongs to a ball centered on  $T_j$  with radius  $\epsilon$ . Alternatively, we could also define  $T_i$  as a nearest neighbour of  $T_j$  if it is one of the  $K$  closest points (measured by the euclidean distance) to  $T_j$  in the set. ii) After the nearest neighbours are defined, they are connected via weighted edges where the weight is the euclidean distance between connected points. iii) Then, the minimum graph distance between each pair of points is computed. This distance is used as a fair approximation to the true geodesic distance (for discussion and proofs, see Tenenbaum et al. (2000) and references therein). The crucial point lies on finding an interval of  $\epsilon$ , or a number of  $K$ , where the solution is robust. Low values of  $\epsilon$  or  $K$  will not connect all the points, while too many will overestimate the dimension of the manifold. After the new matrix of squared distances is computed, the MDS procedure is applied, starting at Eq. (2). The dimensionality of the manifold (the optimum number of dimensions needed to capture the variability of the data) can be measured via the eigenvalues of the MDS procedure. These eigenvalues are a measure of the error made when we project the whole dataset onto the directions defined by the corresponding eigenvectors. The cumulative variance  $c_p$ , of dimension  $p$  is defined using  $c_p = \frac{\sum_{i=1}^p \Lambda_{ii}}{\operatorname{Tr}(\Lambda)}$ , and takes the value  $c_p = 0$  if no statistical variance is explained, and  $c_p = 1$  if all the variability is taken into account.

The Isomap algorithm has two computational bottlenecks (Silva and Tenenbaum, 2003). The first is the calculation of the  $n^2$  shortest-paths distance matrix. The second is the eigenvalue calculation after the double-centering operation in Eq. (2) of the  $n^2$  rank matrix  $\mathbf{Z}^{(2)}$ . These two inefficiencies can be avoided by designating  $n' < n$  landmark points.



**Fig. 2.** The cumulative variances  $r_p$  plotted against the number of dimensions  $p$  considered for PCA (triangles) and Isomap (circles). For (a) spiral, the dimensionality is one when Isomap is used, and two when PCA is applied. For (b) KE SST database, the dimensionality shown by Isomap is three, as the cumulative variance ceases to increase significantly if adding more dimensions. The dimensionality shown by PCA is unclear because of the slow convergence of  $c_p$ .

Instead of computing the whole set of distances, only the  $n \times n'$  matrix of distances from each data point to the landmarks are calculated. Of course, if  $n' \ll n$ , a lot of computing time is saved. The fact that the use of landmarks is feasible can be justified from the assumption that the data are embedded in a low-dimensional manifold (Silva and Tenenbaum, 2003).

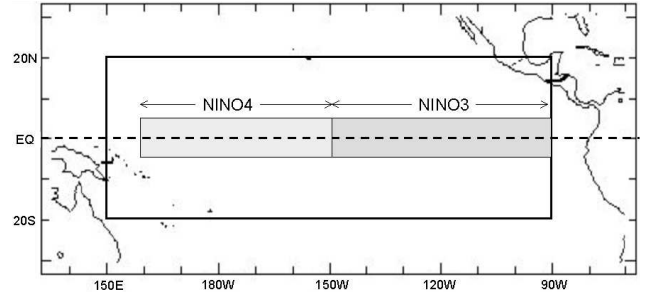
Using the same example as in the preceeding section, we apply Isomap ( $k = 5$ ) to the spiral set with noise represented in Fig. 1a. For the spiral, the dimensionality is one (Fig. 2a), in contrast with the results offered by PCA. We can observe how the spiral is unfolded into an approximately one dimensional set in the Fig. 1c.

In the following section, we will apply PCA and Isomap to geophysical data. More specifically, the SST data in the tropical Pacific Ocean region depicted in Fig. 3 will be analysed using both methods.

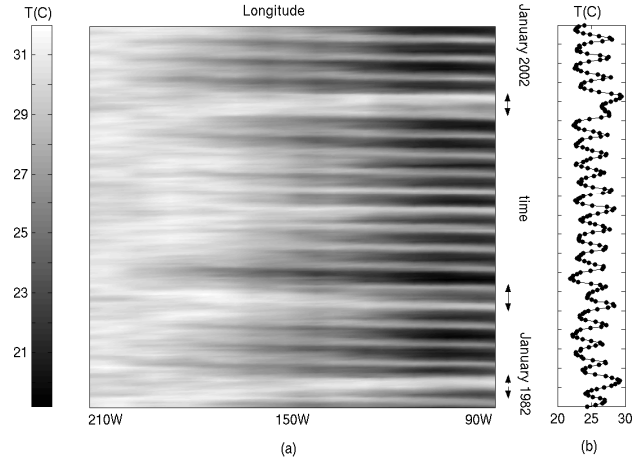
### 3 Application to SST

We have taken Sea Surface Temperature (SST) data from public databases (e.g. <http://ingrid.ldeo.columbia.edu/>). More precisely, we have made use of two databases:

1. the Reynolds-Smith (shortly, RS) database (Reynolds and Smith, 1995) for the region limited by  $89.5^\circ$  W to  $149.5^\circ$  E and  $20.5^\circ$  N to  $20.5^\circ$  S, from November 1981 to October 2002, with a resolution of one month in time and one degree in space. This means that there are 5124 spatial measurements for 252 months.
2. the Kaplan Extended (shortly, KE) database (Kaplan et al., 1998) for the region limited by  $87.5^\circ$  W to  $147.5^\circ$  E and  $17.5^\circ$  N to  $17.5^\circ$  S, from January 1856 to October 1981, with a resolution of one month in time and five degrees in space. In this case, the spatial resolution is 208 points with 1762 time points.

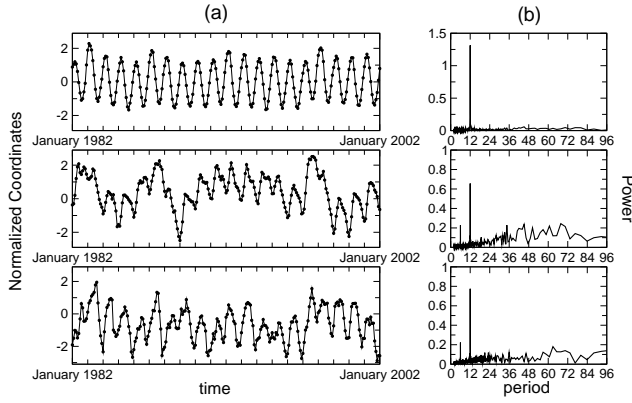


**Fig. 3.** Map of the Pacific Ocean. The big rectangle shows the region under study in comparison to the regions where NINO3 and NINO4 are defined in the literature.

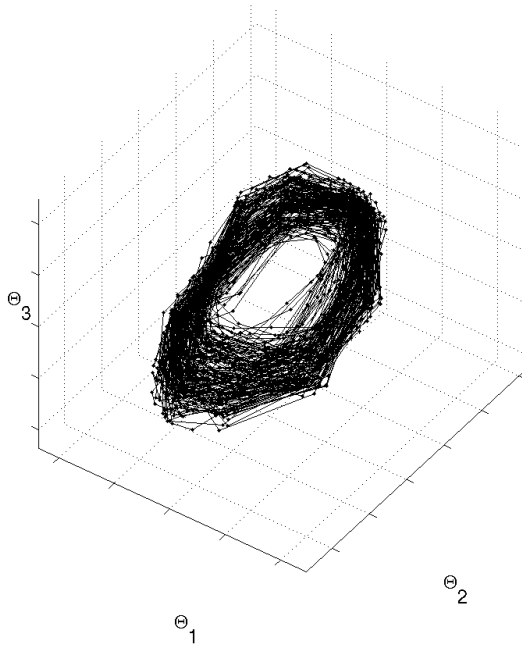


**Fig. 4.** (a) Evolution of the temperature pattern along the Equator from January 1982 to January 2002 (resolution one month). We can observe the warm pool in the west and the cold tongue in the east punctuated by ENSO events (marked with arrows). (b) Time evolution of the temperature at the point  $(115^\circ$  E,  $0^\circ$  N).

RS is used to complement the data of KE in the November 1981 to October 2002 time range. The time evolution of the temperature on the Equator is shown in Fig. 4a. If we focus on a fixed position in the ocean, we can observe two principal oscillations, a rather regular one associated with the annual cycle, and an irregular one associated with ENSO (Fig. 4b). The following results are obtained by using the normalized (unit variance and zero mean) KE database. This allows us to compare PCA and Isomap using the same normalised data. The results obtained from the analysis of the RS database are consistent with the ones presented here. The analysis of SST with PCA shows an annual oscillation which is present in all the principal components (Fig. 5). The convergence of  $c_p$  with the number of dimensions is slow. For this reason, it is difficult to select the number of dimensions that best describes the physical process. Depending of the criterion selected, we get different cut-offs in the number of relevant components. We turn now our efforts to Isomap to compare its results in terms of dimensionality. For the KE database, the radius of the ball that defines the nearest neighbours was taken as  $\epsilon = 9$ . Similar results were found when



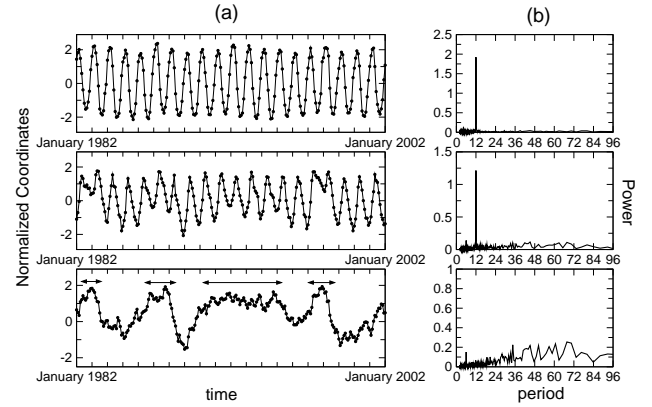
**Fig. 5.** (a) Coordinates and (b) power spectra of, from top to bottom, the first, second and third principal components for the KE dataset from January 1982 to January 2002 (resolution one month) computed via PCA.



**Fig. 6.** Embedding of the KE SST into the dimensions defined by the first three eigenvalues.

$\epsilon$  was ranging from 9 to 30. Although the use of landmark points is suggested for  $n > 1000$ , we used the whole set of data points for calculating the eigenvalues of  $\mathbf{Z}^{(2)}$ . This was done so because, for the KE database,  $n = 1762$ , and the use of the Floyd's algorithm (Floyd, 1962) when computing distances saved enough time of computation without sacrificing accurateness in the calculation of the eigenvalues of  $\mathbf{Z}^{(2)}$ , as they are needed for finding the dimensionality of the data.

The cumulative variances calculated by PCA and Isomap are shown in Fig. 2b. We observe that the dimensionality found by Isomap is three, while the convergence of PCA's cumulative variance is much slower. The first three components found by Isomap are:



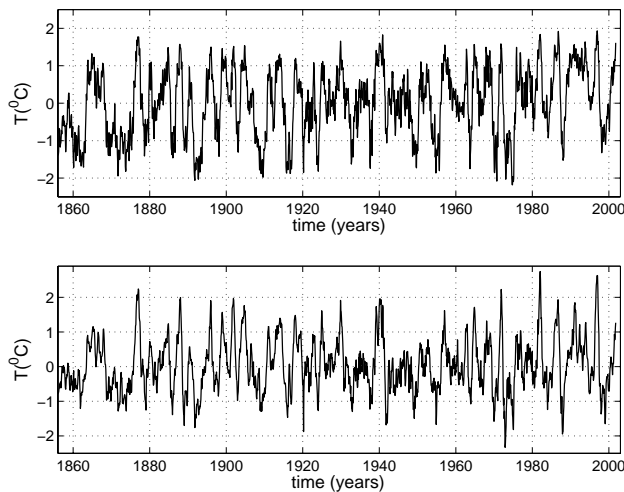
**Fig. 7.** (a) Coordinates and (b) power spectra of, from top to bottom, the first, second and third principal components for the KE dataset from January 1982 to January 2002 (resolution one month) computed via Isomap. The double headed arrows indicate the approximate duration of El Niño events.

$$\begin{aligned} \mathcal{T}_1(t_j) &= \lambda_1^{1/2} \mathbf{p}_{1j} \\ \mathcal{T}_2(t_j) &= \lambda_2^{1/2} \mathbf{p}_{2j} \\ \mathcal{T}_3(t_j) &= \lambda_3^{1/2} \mathbf{p}_{3j}, \end{aligned} \quad (11)$$

which we can plot in a three dimensional coordinate space (Fig. 6). We can observe a twelve-month oscillation in a plane with some deviation in a coordinate perpendicular to that plane. To isolate the twelve-month oscillation we can rotate the new three dimensional points because the representation in terms of distances is equivalent to any orthonormal transformations. The optimal plane can be found by computing the plane that best embeds a mean 12-month cycle over the whole trajectory. After rotating, the three time series shown in Fig. 7 are extracted. We can see that there are major differences between the second and third components of Figs. 5 and 7. Moreover, the third Isomap component faithfully represents ENSO, as all the well known events have their corresponding peak in the time series of the third component.

The results show that the complex dynamics due to the interaction between ENSO and the annual cycles can be well approximated by a three dimensional manifold.

PCA introduces a particular normalisation in the data because the matrix of the eigendecomposition is the correlation, which is naturally normalised to variance one. But we can apply MDS and Isomap to the raw data, without any other mathematical operation. The results are similar to the usual indexes that describe ENSO, as we can see in Fig. 8. This plot shows how this decomposition provides a useful way to characterise ENSO by using the third Isomap component. Moreover, we observe that the occurrence of the events is essentially preserved, although the amplitude and the probability distribution function found by the decomposition are slightly different. This is due to the fact that now the annual cycle is not approximated by a periodic function, as in the NINO3.4 or other indexes.



**Fig. 8.** Third Isomap coordinate (top) in comparison with the index NINO3.4 (bottom). The amplitude of the events changes slightly from one method to the other.

#### 4 Conclusions

Nonlinear dimensionality reduction methods provide a useful way of analysing and modeling high dimensional data when nonlinear interactions are present. If the physical process can be embedded in a low dimensional manifold, the reduction of the relevant components is better achieved by nonlinear methods than by linear ones. In particular, Isomap provides a physically appealing method of decomposing the data, as it substitutes the euclidean distances in the manifold by an approximation of the geodesic distances. We expect that this method could be successfully applied to other oscillatory extended systems and, in particular, to meteorological phenomena. Reduction allows to model the physical process by a low dimensional nonlinear dynamical system strictly based on data in order to make predictions of future states. We expect that this model will provide insight into the physics of the interaction between ENSO and the annual cycle in the tropical Pacific Ocean.

**Acknowledgements.** We thank U. Parlitz for guiding us to the Isomap algorithm. We also thank J. Wichard for discussions. A. J. G. and J. K. were supported by the European Training Network COSYC of SENS (Contract Number HPRN-CT-2000-00158). A. T. acknowledges support from the Deutsche Forschungsgemeinschaft through the Collaborative Research Effort SFB460.

Edited by: M. Thiel

Reviewed by: A. Chian and another referee

#### References

An, S.-I., and Wang, B.: Mechanisms of locking of the El Niño and La Niña mature phases to boreal winter, *J. Climate*, 14, 2164–2176, 2001.

Borg, I. and Groenen, P.: *Modern Multidimensional Scaling: Theory and Applications*, Springer-Verlag, 1997.

Do Carmo, M.: *Differential Geometry of Curves and Surfaces*, Prentice-Hall, 1976.

Floyd, R. W.: Algorithm 97: Shortest path, *Communications of the ACM*, 5, 6, 345, 1962.

Gower, J.: Some distance properties of latent root and vector methods used in multivariate analysis, *Biometrika*, 53, 325–328, 1966.

Grieger, B. and Latif, M.: Reconstruction of the El Niño attractor with neural networks, *Climate Dynamics*, 10, 267–276, 1994.

Holmes, P., Lumley, J. L., Berkooz, G., Mattingly, J. C., and Wittenberg, R. W.: Low-dimensional models of coherent structures in turbulence, *Physics Reports*, 287(4), 337–384, 1997.

Jin, F. F.: El Niño/Southern Oscillation and the annual cycle: sub-harmonic frequency-locking and aperiodicity, *Physica D*, 98, 442–465, 1996.

Jin, F. F.: An equatorial recharge paradigm for ENSO. Part I: Conceptual model, *J. Atmos. Sci.*, 54, 811–829, 1997.

Jolliffe, I. T. *Principal Component Analysis*, Springer Verlag, 1986.

Kaplan, A., Cane M., Kushnir, Y., Clement, A., Blumenthal, M., and Rajagopalan, B.: Analyses of global sea surface temperature 1856–1991, *J. Geophys. Res.*, 103(C9), 18, 567–18, 589, 1998.

Karhunen, K.: Zur Spektraltheorie stochastischer Prozesse, *Ann. Acad. Sci. Fennicae, Ser. A1*, 34, 1946.

Li, T. and Philander, S. G.: On the Annual Cycle of the Eastern Equatorial Pacific, *J. Climate*, 9, 2986–2998, 1996.

Loève, M.: *Fonctions aléatoires de second ordre*, *Comptes Rendus Acad. Sci., Paris*, 220, 1945.

Monahan, A. H.: Nonlinear principal component analysis: Tropical Indo-Pacific sea surface temperature and sea level pressure, *J. Climate*, 14, 219–233, 2001.

Philander, S. G.: *El Niño, La Niña and the Southern Oscillation*, San Diego: Academic Press, 1990.

Pikovsky, A., Rosenblum, M., and Kurths, J.: *Synchronization: A universal concept in nonlinear sciences*, Cambridge University Press, 2001.

Reynolds, R. W. and Smith, T. M.: A high resolution global sea surface temperature climatology, *J. Climate*, 8, 1571–1583, 1995.

von Storch, H. and Zwiers, F. W.: *Statistical Analysis in Climate Research*, Cambridge University Press, 1999.

Tenenbaum, J. B., de Silva, V., and Langford, J. C.: A global geometric framework for nonlinear dimensionality reduction, *Science*, 290, 2319–2323, 2000.

Tziperman, E., Stone, L., Cane, M., and Jarosh, H.: El Niño chaos: overlapping of resonances between the seasonal cycle and the Pacific ocean-atmosphere oscillator, *Science*, 264, 72–74, 1994.

Tziperman, E., Cane, M. A., Zebiak, S. E., Xue, Y., and Blumenthal, B.: Locking of El Niño's Peak Time to the End of the Calendar Year in the Delayed Oscillator Picture of ENSO, *J. Climate*, 11, 2191–2199, 1998.

de Silva, V. and Tenenbaum, J. B.: Global versus local methods in nonlinear dimensionality reduction, *Advances in Neural Information Processing Systems*, edited by Becker, S., Thrun, S., and Obermayer, K., MIT Press, Cambridge, 15, 705–712, 2003.

Xie, S.: On the Genesis of the Equatorial Annual Cycle, *J. Climate*, 7, 2008–2013, 1994.

Xie, S.: Interaction between the Annual and Interannual Variations in the Equatorial Pacific, *J. Phys. Oceanogr.*, 25, 1930–1941, 1995.

Zebiak, S. E. and Cane, M. A.: A model El Niño/Southern Oscillation, *Mon. Weather Rev.*, 115, 2262–2278, 1987.